

Lecture Notes in
Statistical and Mathematical Methods for
Artificial Intelligence
91255

Elena Loli Piccolomini

elena.loli@unibo.it

Academic Year 2020/2021

Last updated: September 2020

Thanks to: Giorgio Renzi

Lorenzo Mario Amorosa

Lucia La Forgia

Giacomo Pinardi

for the help provided

Contents

1	Linear systems	1
1.1	Introduction to linear systems	1
1.2	Square linear systems	1
1.3	Direct methods	2
1.4	Iterative methods	3
1.5	Inherent errors in linear systems	5
1.6	Linear least squares	6
1.6.1	Orthogonality and projections	6
1.6.2	The linear least squares problem	8

1 Linear systems

1.1 Introduction to linear systems

A linear system can be written as

$$\mathbf{A}\mathbf{x} = \mathbf{b}$$

where \mathbf{A} is a matrix of size $m \times n$ (we suppose $m \geq n$) and \mathbf{x} is a column vector of length n and \mathbf{b} is a column vector of length m . \mathbf{x} represents the unknown solution while \mathbf{b} is a given vector. This form can be expanded as

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n &= b_1 \\ a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n &= b_2 \\ &\vdots \\ a_{m1}x_1 + a_{m2}x_2 + \dots + a_{mn}x_n &= b_m \end{aligned}$$

We are interested in:

- Existence and uniqueness of a solution
- Numerical methods to find the solution
- Conditioning of the problem

We separately consider the case of square linear systems ($m = n$) and least squares problem ($m > n$).

1.2 Square linear systems

Proposition 1.1. The solution of a linear system

$$\mathbf{A}\mathbf{x} = \mathbf{b}$$

with \mathbf{A} of dimension $n \times n$, \mathbf{x} and \mathbf{b} of dimension n exists and is unique iff one of the following conditions holds:

- \mathbf{A} is non-singular
- $\text{rank}(\mathbf{A}) = n$
- The system $\mathbf{A}\mathbf{x} = \mathbf{0}$ admits only the solution $\mathbf{x} = \mathbf{0}$

The solution can be algebraically computed in the following way

$$\mathbf{Ax} = \mathbf{b} \implies \mathbf{A}^{-1}\mathbf{Ax} = \mathbf{A}^{-1}\mathbf{b} \implies \mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$$

The problem with this kind of method is that computing the inverse of \mathbf{A} can be expensive for large values of n .

Another method to compute the solution of a system is *Cramer's rule*:

$$x_i = \frac{\det A_i}{\det A}, \quad i = 1, \dots, n$$

where A_i is obtained from \mathbf{A} by replacing the i -th column by \mathbf{b}

As for the previous one, the computation of determinants can be quite expensive for large values of n .

We will consider two different approaches to find the solution of a linear system:

- *direct methods*. They yield the solution in a finite number of steps; they are more precise but more expensive in terms of computational cost
- *iterative methods*. They require (in principle) an infinite number of steps; they are less precise and less expensive

It is very important to perform also an error analysis on the obtained results. It generally includes two parts:

- *Rounding or arithmetic errors*. They depend on the algorithm steps. An algorithm producing an arithmetic error limited by a constant is called a **stable algorithm**.
- *Inherent errors*. They are due to the errors in the data representation and the DO NOT depend on the algorithm. For this reason, if the inherent error is large, we speak about **ill-posed problem**.

1.3 Direct methods

In the case of direct methods we factorize the matrix \mathbf{A} as in section ?? . Let $\mathbf{A} = \mathbf{LU}$, then the solution of the linear system $\mathbf{Ax} = \mathbf{b}$ can be computed by solving two triangular systems

$$\begin{aligned} \mathbf{Ly} &= \mathbf{b} \text{ (forward substitutions algorithm)} \\ \mathbf{Ux} &= \mathbf{y} \text{ (backward substitutions algorithm)} \end{aligned}$$

Forward substitution.

Esempio 2.1 *Sostituzione in avanti.*

$$\begin{pmatrix} \ell_{1,1} & 0 & 0 \\ \ell_{2,1} & \ell_{2,2} & 0 \\ \ell_{3,1} & \ell_{3,2} & \ell_{3,3} \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \\ y_3 \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \\ b_3 \end{pmatrix}$$

$$\begin{array}{ll} y_1 = b_1 / \ell_{1,1} & 1 \text{ operazione} \\ \downarrow & \\ y_2 = (b_2 - \ell_{2,1} * y_1) / \ell_{2,2} & 3 \text{ operazioni} \\ \downarrow & \\ y_3 = (b_3 - \ell_{3,1} * y_1 - \ell_{3,2} * y_2) / \ell_{3,3} & 5 \text{ operazioni} \end{array}$$

Backward substitution.

$$\begin{pmatrix} u_{1,1} & u_{1,2} & u_{1,3} \\ 0 & u_{2,2} & u_{2,3} \\ 0 & 0 & u_{3,3} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} y_1 \\ y_2 \\ y_3 \end{pmatrix}$$

$$\begin{array}{ll} x_3 = y_3 / u_{3,3} & 1 \text{ operazione} \\ \downarrow & \\ x_2 = (y_2 - u_{2,3} * x_3) / u_{2,2} & 3 \text{ operazioni} \\ \downarrow & \\ x_1 = (y_1 - u_{1,2} * x_2 - u_{1,3} * x_3) / u_{1,1} & 5 \text{ operazioni} \end{array}$$

In the case of pivoting algorithm, we have that $PA = LU$. Then the solution of the linear system $PA\mathbf{x} = P\mathbf{b}$ can be computed by solving two triangular systems

$$\begin{array}{l} L\mathbf{y} = P\mathbf{b} \text{ (forward substitutions algorithm)} \\ U\mathbf{x} = \mathbf{y} \text{ (backward substitutions algorithm)} \end{array}$$

1.4 Iterative methods

Iterative methods are numerical methods that approximate a solution using a procedure involving several (or an infinite amount of) steps.

The basic idea of iterative methods is to construct a sequence of vectors \mathbf{x}_k enjoying the property of *convergence*

$$\mathbf{x}^* = \lim_{k \rightarrow \infty} \mathbf{x}_k$$

where \mathbf{x}^* is the exact solution and the starting guess \mathbf{x}_0 is given.

In general the sequence \mathbf{x}_k is obtained as $\mathbf{x}_k = g(\mathbf{x}_{k-1})$, $k = 1, \dots$ where g is a particular function or set of operations acting on \mathbf{x}_{k-1} . Different classes of iterative methods are defined for different expressions of g . In general, an iterative method converges for matrices with fixed properties (such as on the shape, on the eigenvalues, ...).

Iterative methods are sensitive to both algorithmic and truncation errors.

Examples of iterative methods (the most common):

- *stationary iterative methods*, they take the form

$$\mathbf{x}_{k+1} = \mathbf{B}\mathbf{x}_k + \mathbf{f}$$

where \mathbf{B} is called *iteration matrix* and \mathbf{f} is a vector obtained from \mathbf{b}

- *gradient-like methods*, they take the form

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{p}_k$$

where $\alpha_k \in \mathbb{R}$, \mathbf{p}_k is a vector called *direction*.

If \mathbf{A} is symmetric and positive definite and the vectors \mathbf{p}_k have the *conjugacy* property, i.e:

$$\mathbf{p}_i^T \mathbf{A} \mathbf{p}_k = 0 \quad \text{if } i \neq k$$

then the method is called *Conjugate Gradients*.

All these methods require one matrix-vector multiplication per iteration. Hence, if \mathbf{A} is a *full* matrix the computational complexity is $O(n^2)$ per iteration.

Stopping criteria for iterative methods for linear systems are:

- defined the residual $\mathbf{r}_k = \mathbf{b} - \mathbf{A}\mathbf{x}_k$ at iteration k , stop the algorithm when
 - $\|\mathbf{r}_k\| \leq \epsilon$ (absolute criterion)
 - $\frac{\|\mathbf{r}_k\|}{\|\mathbf{b}\|} \leq \epsilon$ (relative criterion)
- $\|\mathbf{x}_{k+1} - \mathbf{x}_k\| \leq \tau$ (absolute criterion)

Sparse matrices A matrix A is called *sparse* if it has a very small percentage (i.e. $0.05 - 0.06\%$) of non-null elements. In this case, the matrix is not fully stored, but only the non-null elements and their indices are stored. The matrix-vector multiplication requires k multiplications, where k is the number of non-null elements.

1.5 Inherent errors in linear systems

We remind that inherent errors are due to the errors in the data representation. They **do not** depend on the algorithm used for computing the result. Really, the analysis of inherent errors is performed by supposing to use exact arithmetic to compute the result.

Consider now a linear system $A\mathbf{x} = \mathbf{b}$. What happens to the solution \mathbf{x} when A and/or \mathbf{b} slightly change (e.g. due to machine precision)?

Suppose that A doesn't change while \mathbf{b} changes in $\mathbf{b} + \Delta\mathbf{b}$. The linear system becomes

$$A(\mathbf{x} + \Delta\mathbf{x}) = \mathbf{b} + \Delta\mathbf{b}, \Delta\mathbf{x} \text{ inherent error}$$

We want to compare $\left\| \frac{\Delta\mathbf{x}}{\mathbf{x}} \right\|$ and $\left\| \frac{\Delta\mathbf{b}}{\mathbf{b}} \right\|$ to see how much the solution has changed with respect to \mathbf{b} . Let's subtract $A\mathbf{x} = \mathbf{b}$ to $A(\mathbf{x} + \Delta\mathbf{x}) = \mathbf{b} + \Delta\mathbf{b}$

$$A\Delta\mathbf{x} = \Delta\mathbf{b} \implies \Delta\mathbf{x} = A^{-1}\Delta\mathbf{b}$$

Then we have

$$\|\Delta\mathbf{x}\| = \|A^{-1}\Delta\mathbf{b}\| \leq \|A^{-1}\| \|\Delta\mathbf{b}\|$$

and from the linear system equation

$$\|A\mathbf{x}\| = \|\mathbf{b}\| \implies \|\mathbf{b}\| = \|A\mathbf{x}\| \leq \|A\| \|\mathbf{x}\| \implies \|\mathbf{x}\| \geq \left\| \frac{\mathbf{b}}{A} \right\|$$

Thus we have

$$\left\| \frac{\Delta\mathbf{x}}{\mathbf{x}} \right\| \leq \|A^{-1}\| \cdot \|A\| \cdot \frac{\|\Delta\mathbf{b}\|}{\|\mathbf{b}\|} = K(A) \frac{\|\Delta\mathbf{b}\|}{\|\mathbf{b}\|}$$

Definition 1.1. The *condition number* of a matrix $A \in \mathbb{C}^{n \times n}$ is defined as

$$K(A) = \|A\| \|A^{-1}\|$$

where $\|\cdot\|$ is a matrix norm. In general $K(A)$ depends on the choice of the norm, indicated by a subscript.

Notice that $K(A) \geq 1$ since

$$1 = \|AA^{-1}\| \leq \|A\| \|A^{-1}\| = K(A)$$

The condition number measures how sensitive a function is to changes/errors in the input and how much the output changes as a result of the operations performed. A *well-conditioned* system is a system where $K(A)$ is small, while an *ill-conditioned* system is a system where $K(A)$ is large.

If $K(A)$ is very large, the matrix A is near (with respect to a norm) a singular matrix and its columns are quasi-linearly dependent. Regularization techniques can be used in order to reduce $K(A)$.

Example 1.1. The condition number for p-norm with $p = 2$ is

$$\begin{aligned} \|A\|_2 &= \sqrt{\rho(A^T A)} = \sigma_{max} \\ \|A^{-1}\|_2 &= \frac{1}{\sqrt{\rho(A^T A)}} = \frac{1}{\sigma_{min}} \\ K_2(A) &= \|A\| \|A^{-1}\| = \frac{\sigma_{max}}{\sigma_{min}} \end{aligned}$$

where λ_{max} and λ_{min} are the maximum and minimum singular values of A

1.6 Linear least squares

Consider an *overdetermined system*

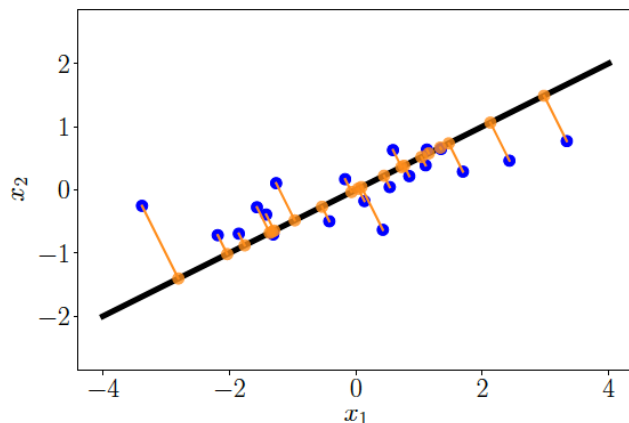
$$A\mathbf{x} = \mathbf{b}$$

where A has dimension $m \times n$ with $m > n$, \mathbf{x} of dimension n and \mathbf{b} of dimension m . Such a system usually has no solution, meaning that \mathbf{b} does not lie on the subspace spanned by the columns of A (denoted by $\text{Col}(A)$ from now on). Since no exact solution exists for this problem, we would like to compute the best approximate solution. We need some tools.

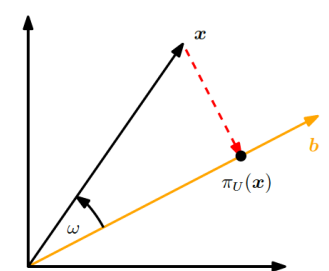
1.6.1 Orthogonality and projections

Projections are an important class of linear transformations. In machine learning one often deals with high-dimensional data which is hard to visualize

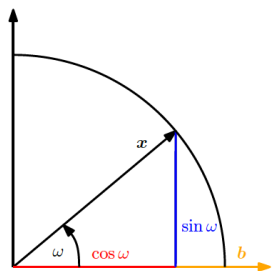
and analyze. Oftentimes, only a few dimensions contain useful information, meaning that other dimensions are not essential to understanding the data. When applying data compression techniques we want to minimize the loss of information by finding the most informative dimensions of the data first.



Orthogonal projection onto a subspace of dimension 1.



(a) Projection of $x \in \mathbb{R}^2$ onto a subspace U with basis vector b .



(b) Projection of a two-dimensional vector x with $\|x\| = 1$ onto a one-dimensional subspace spanned by b .

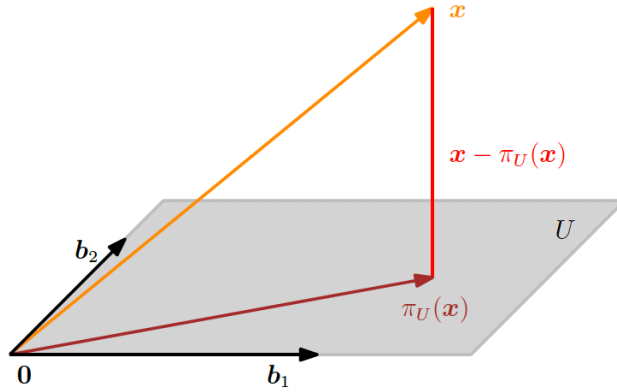
- The projection $z = \Pi_U(\mathbf{x})$ is the closest point to \mathbf{x} of the subspace U generated by \mathbf{b} by considering as the measure of the distance the 2-norm ($\|\cdot\|_2$). The vector $\mathbf{z} - \mathbf{x}$ (the red dashed one in the figure) is orthogonal to the subspace U (and to the vector \mathbf{b} generating U). Hence, the scalar product: $\langle \mathbf{z} - \mathbf{x}, \mathbf{b} \rangle = 0$
- The vector \mathbf{z} is a vector of the subspace generated by \mathbf{b} , hence $\mathbf{z} = \alpha \mathbf{b}$, $\alpha \in \mathbb{R}$.

We can find α in the following way

$$\langle \mathbf{x} - \alpha \mathbf{b}, \mathbf{b} \rangle = 0 \implies \langle \mathbf{x}, \mathbf{b} \rangle - \alpha \langle \mathbf{b}, \mathbf{b} \rangle = 0$$

$$\alpha = \frac{\langle \mathbf{x}, \mathbf{b} \rangle}{\langle \mathbf{b}, \mathbf{b} \rangle} = \frac{\mathbf{x}^T \mathbf{b}}{\|\mathbf{b}\|^2}$$

Example when projecting from \mathbf{R}^3 to \mathbf{R}^2 .



We now generalize to sub-spaces of dimension possibly greater than one.

Theorem 1.1. (*Best approximation theorem*) Let W a subspace with dimension p of \mathbb{R}^n , $\mathbf{y} \in \mathbb{R}^n$ and $\hat{\mathbf{y}} = \text{proj}_W(\mathbf{y})$. Then $\|\mathbf{y} - \hat{\mathbf{y}}\| < \|\mathbf{y} - \mathbf{v}\|$, $\forall \mathbf{v} \in W$, $\mathbf{v} \neq \hat{\mathbf{y}}$.

This means that $\hat{\mathbf{y}}$ is the best approximation of \mathbf{y} in W .

Example 3.10 MML book.

The idea is to use projections to find the vector in $\text{Col}(\mathbf{A})$ (i.e. the subspace generated by set of all the columns of \mathbf{A}) which is closest to \mathbf{b} . As already seen in the previous chapter, this vector is indeed the orthogonal projection of \mathbf{b} onto the subspace $\text{Col}(\mathbf{A})$.

1.6.2 The linear least squares problem

When \mathbf{A} has size $m \times n$ with $m > n$, it is not possible to find a solution of the linear system $\mathbf{A}\mathbf{x} = \mathbf{b}$. The orthogonal projections allow to find a vector \mathbf{x} that is *near* the system solution (in the sense of the 2-norm).

The problems is then formulated as to compute $\hat{\mathbf{x}}$ so that:

$$\hat{\mathbf{x}} = \operatorname{argmin}_x \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2 \quad (1)$$

For what concerns existence and uniqueness of the solution of the linear least squares problem, two different cases are possible:

- $\operatorname{rank}(\mathbf{A}) = n$, meaning that every column of \mathbf{A} is linearly independent. A unique solution exists $\forall \mathbf{b} \in \mathbb{R}^m$. Let \mathbf{x}^* be the approximate solution to the linear system, we have

$$\mathbf{A}\mathbf{x}^* = \operatorname{proj}_{\operatorname{Col}(\mathbf{A})}(\mathbf{b})$$

This means that

$$\mathbf{b} - \mathbf{A}\mathbf{x}^* \in [\operatorname{Col}(\mathbf{A})]^\perp \implies \mathbf{A}^T(\mathbf{b} - \mathbf{A}\mathbf{x}^*) = 0$$

Thus, we obtain the *normal equations*

$$\mathbf{A}^T \mathbf{A} \mathbf{x}^* = \mathbf{A}^T \mathbf{b}$$

which is a linear system in n equations and n unknowns. $\mathbf{A}^T \mathbf{A}$ is called *Gramian matrix* of \mathbf{A} and is non-singular, positive definite and symmetric. The approximate solution can then be obtained

$$\mathbf{x}^* = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{b}$$

- $\operatorname{rank}(\mathbf{A}) = k$, $k < n$, meaning that there is an infinite number of solutions (∞^{n-k}). We can find the least squares solution to the problem by using the pseudoinverse of \mathbf{A} obtained via SVD

$$\mathbf{A}\mathbf{x} = \mathbf{b} \implies \mathbf{A}^+ \mathbf{A} \mathbf{x}^* = \mathbf{A}^+ \mathbf{b} \implies \mathbf{x}^* = \mathbf{A}^+ \mathbf{b} \implies \mathbf{x}^* = \sum_{i=1}^k \frac{\mathbf{u}_i^T \mathbf{b}}{\sigma_i} \mathbf{v}_i$$

It is possible to modify the linear least squares problem (1) by introducing a weight matrix with the aim to give different weights to the components of \mathbf{x} .

$$\hat{\mathbf{x}} = \operatorname{argmin}_x \|\mathbf{W} \mathbf{A} \mathbf{x} - \mathbf{b}\|_2^2 \quad (2)$$

where \mathbf{W} is an invertible matrix. In this case the *weighted normal equations* are:

$$(\mathbf{W}\mathbf{A})^T \mathbf{W}\mathbf{A} \mathbf{x} = (\mathbf{W}\mathbf{A})^T \mathbf{b}$$